

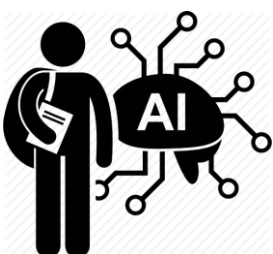


Introduction to NLP for Software Developers

Maciej Szymczak

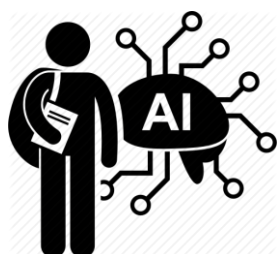
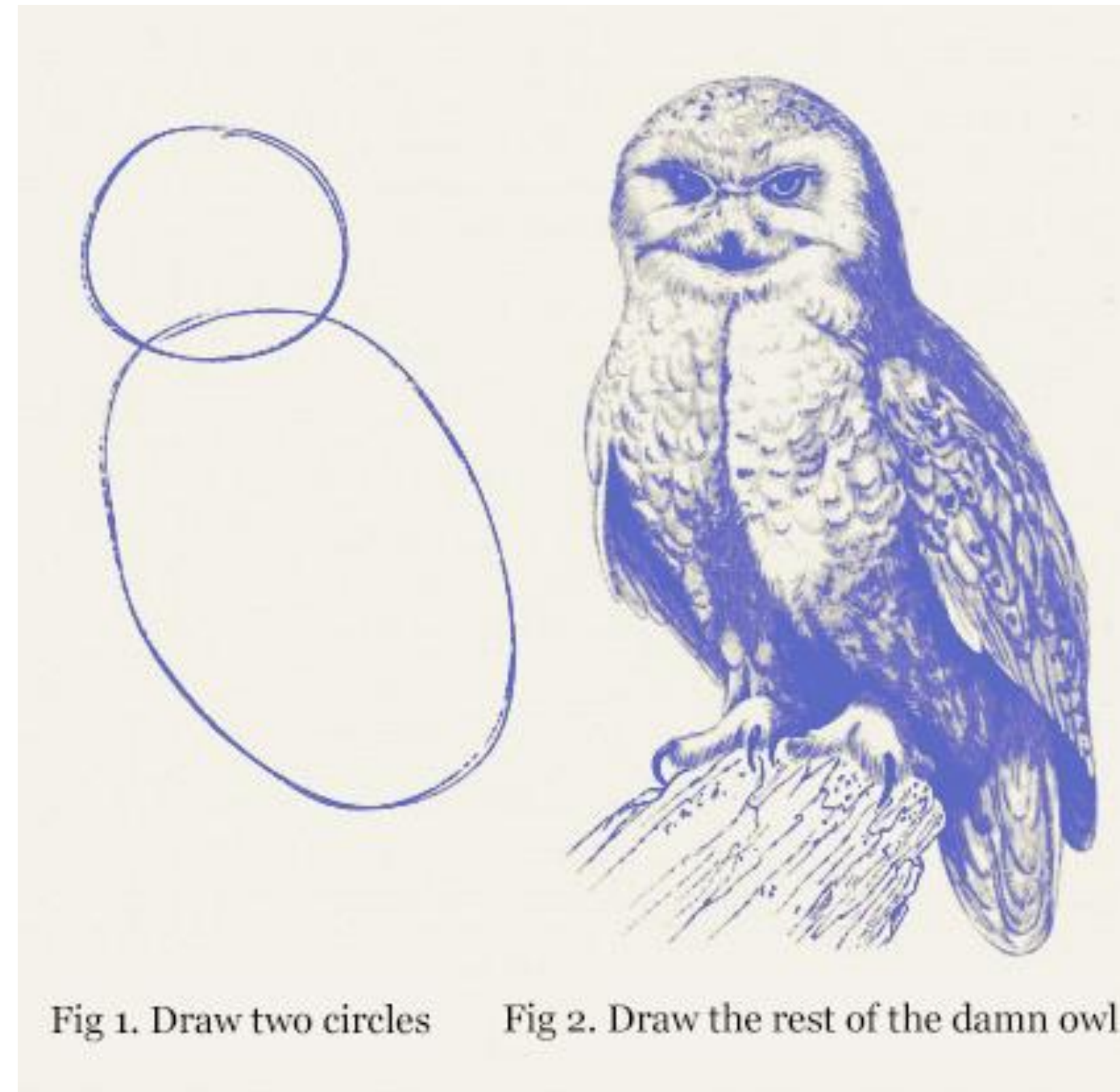
AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP



Agenda

- ✓ What is NLP
- ✓ How to do NLP
- ✓ Some advanced topics
- ✓ DEMO
- ✓ Conclusion and Q&A

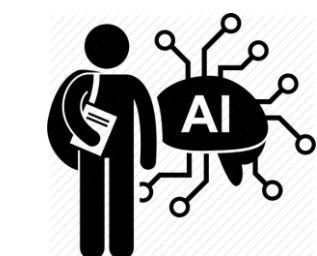
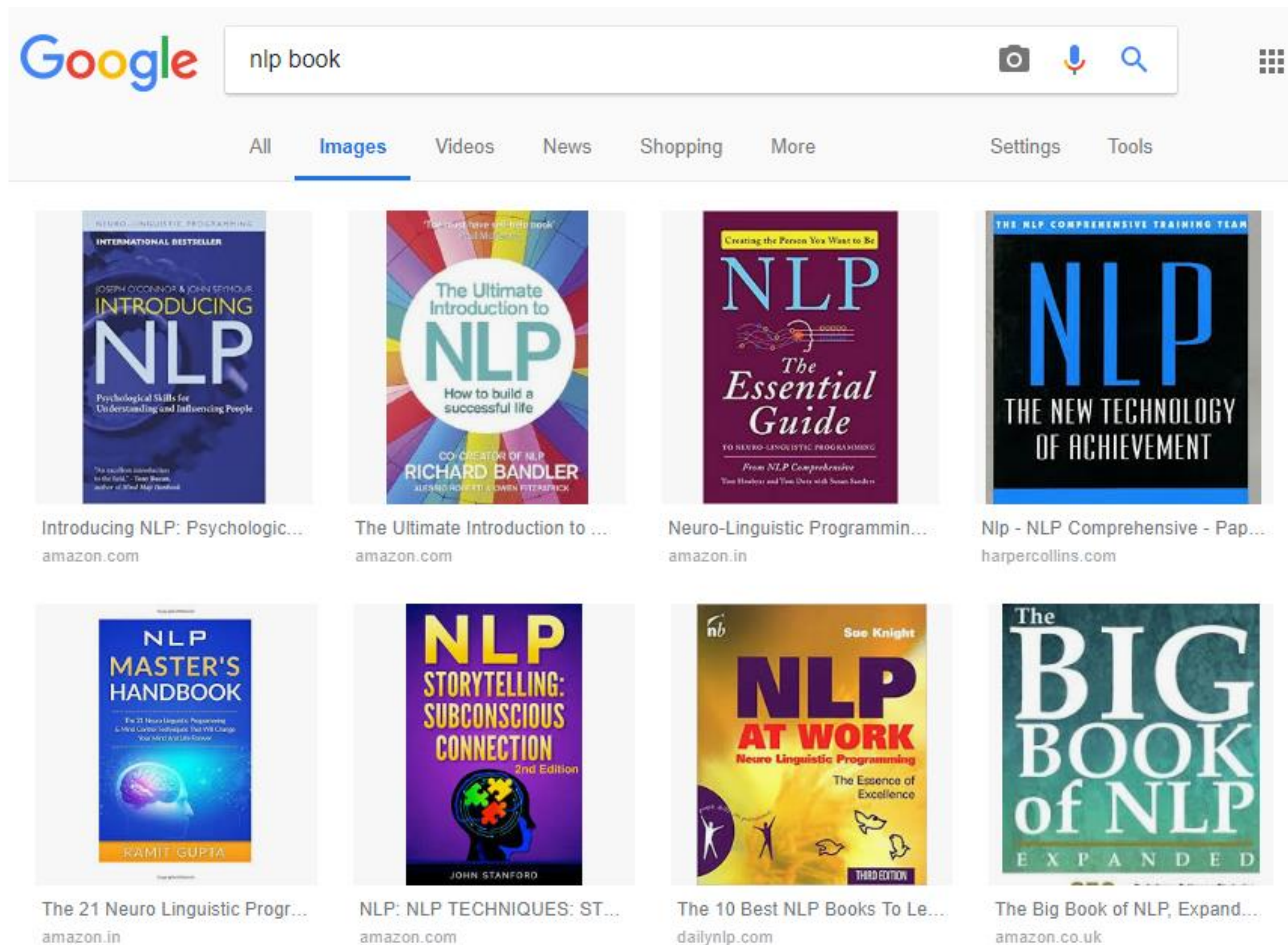




What exactly is NLP?



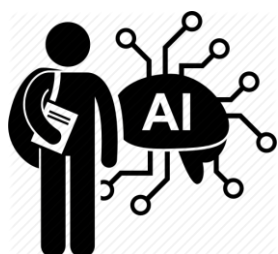
Well... Not this



Rather this...

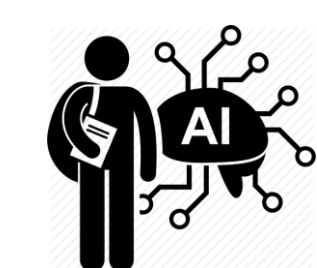


AI Day Wrocław
Volvo Group IT in Poland
29th November 2018





Natural language processing (NLP) is a subfield of computer science (...) concerned with (...) how to program computers to process (...) natural language data.



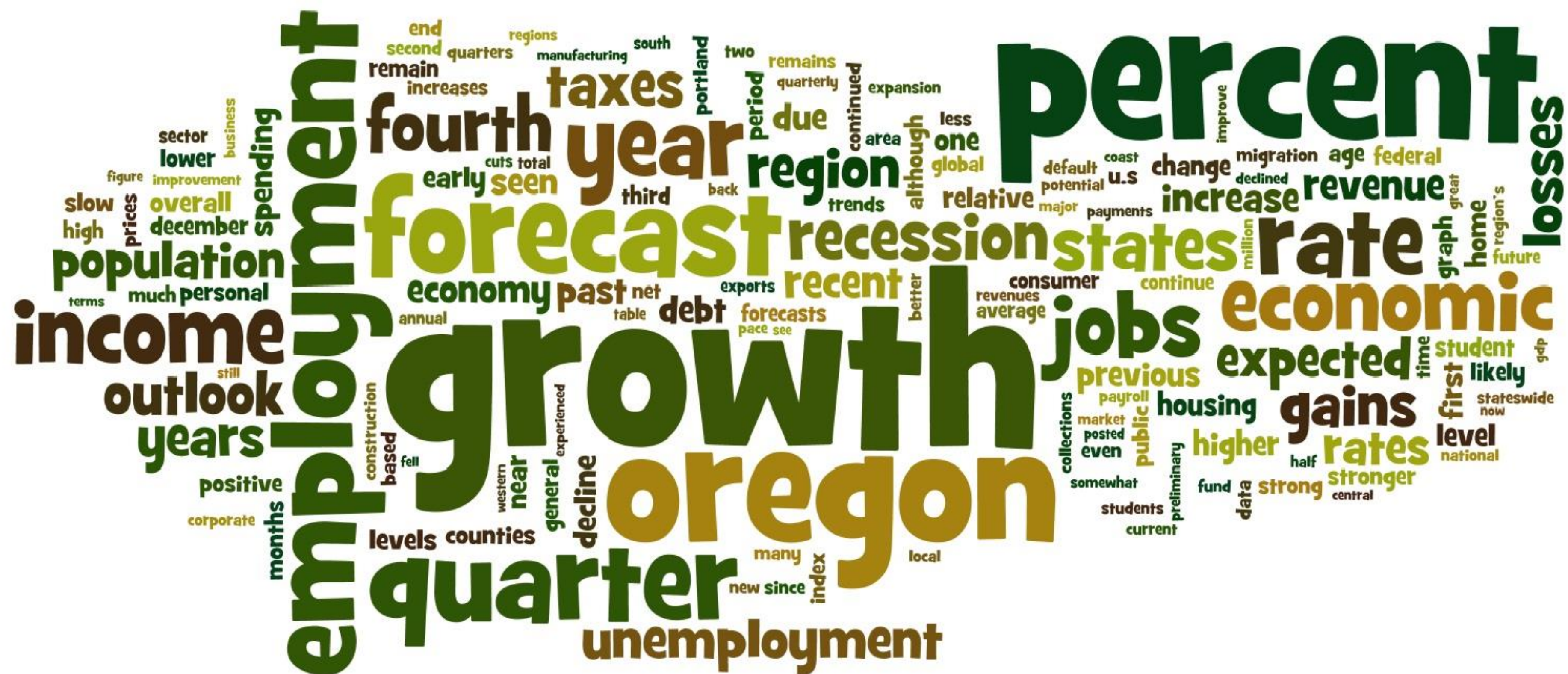


How to achieve that?



AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP



Source: Oregon Office of Economic Analysis

AI Day Wroclaw
Volvo Group IT in Poland
29th November 2018

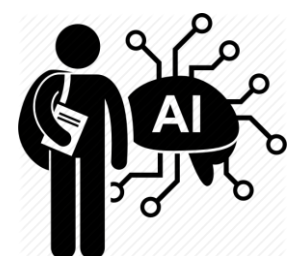
Tokenization



sentence:

"To be or not to be, that is the question."

tokens: | " | To | be | or | not | to | be | , | that | is | the | question | . | " |

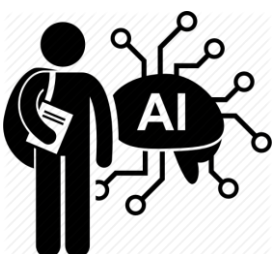


Stemming

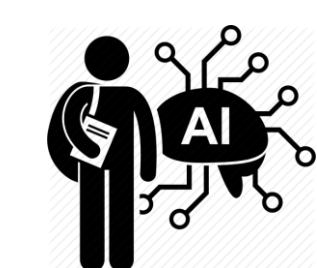


ing \rightarrow ϵ , e.g. running \rightarrow run
ational \rightarrow ate, e.g. rational \rightarrow rate

Sometimes wrong, e.g. *booking!*



Stop words



Bag of Words (BOW)



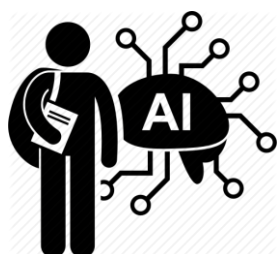
text 1

I had a dream

text 2

Dream a little dream of me

	i	had	a	dream	little	of	me
BoW 1	1	1	1	1	0	0	0
BoW 2	0	0	1	2	1	1	1



Term frequency – inverse document frequency (tf-idf)

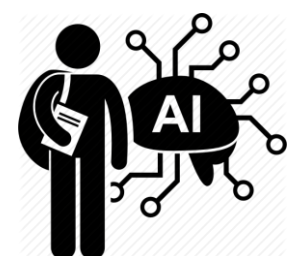


$$tf-idf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

$$idf(t, D) = \ln \left(\frac{|D|}{|\{d \in D : f_{t,d} > 0\}|} \right)$$

$t = \text{term}$, $d = \text{document}$, $D = \text{set of documents}$, $f_{t,d} = \text{count}(t, d)$



N-grams



“What a wonderful world”

| <S> What | What a | a wonderful | wonderful world | world <E>

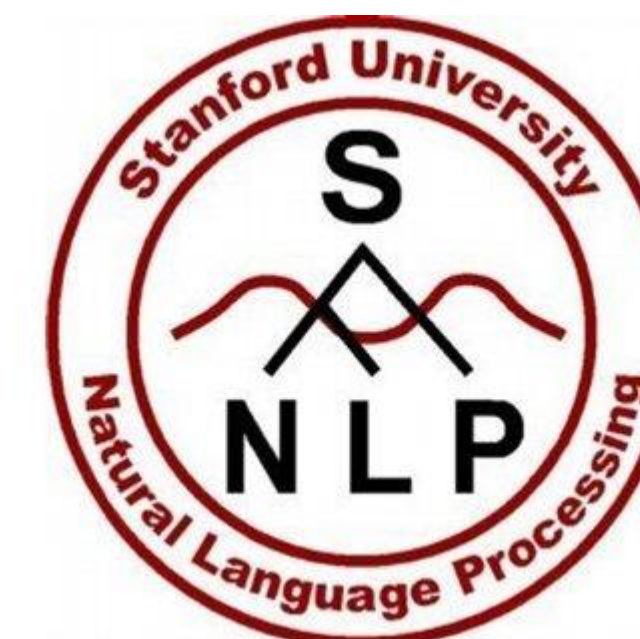


Software packages - code



NLTK

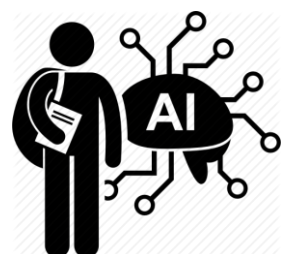
spaCy



gensim



NLP as a service



AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP

Advanced topics



Stemming



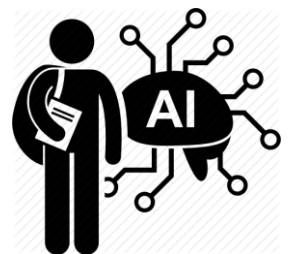
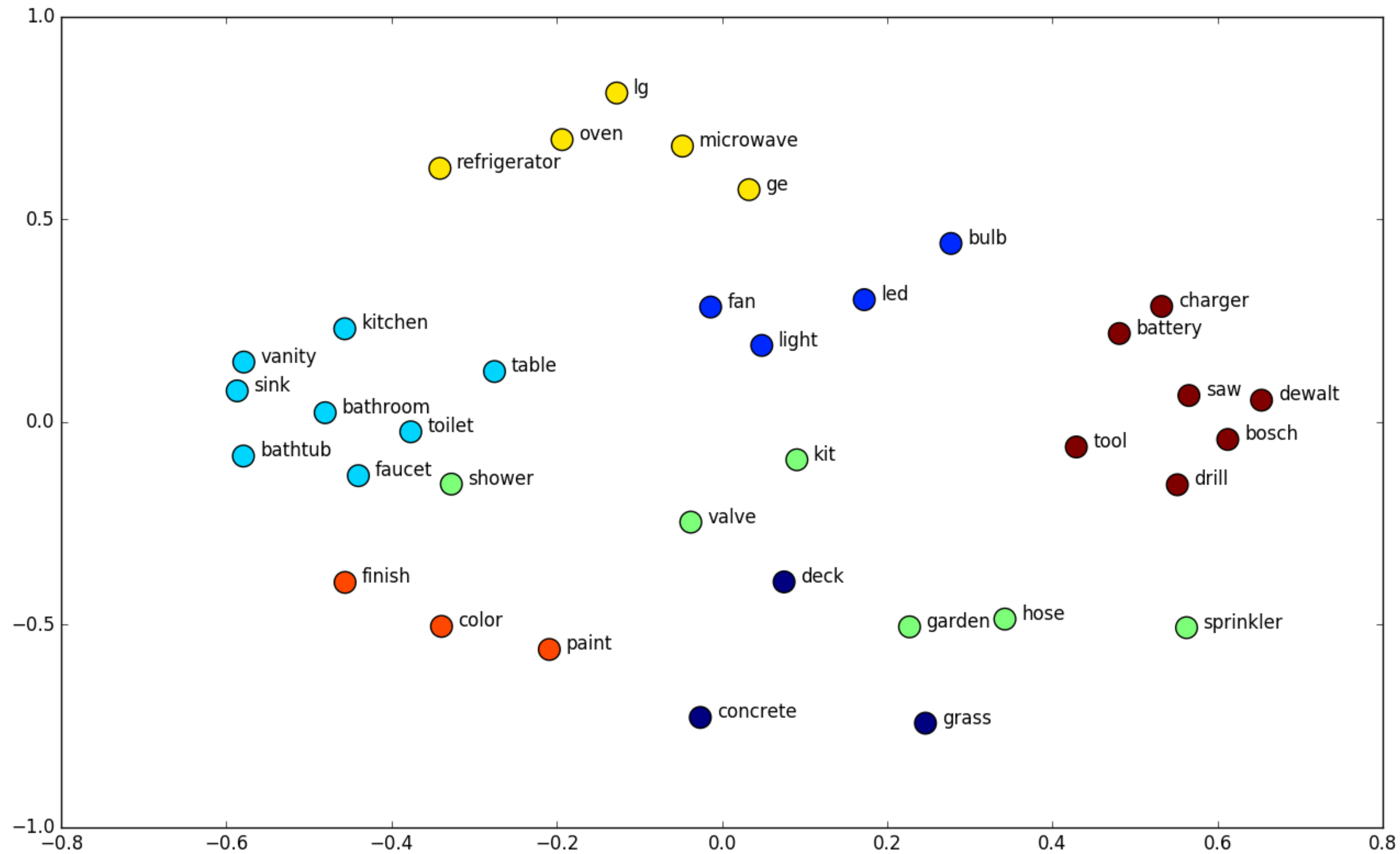
tf-idf



Deep Learning



Word Embeddings



Source: <https://www.shanelynn.ie/get-busy-with-word-embeddings-introduction/>

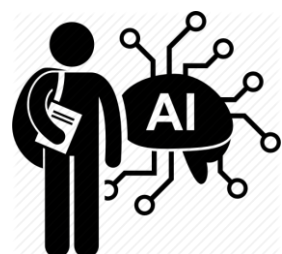
AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP



*„You shall know a word
by the company it keeps”*

John Rupert Firth

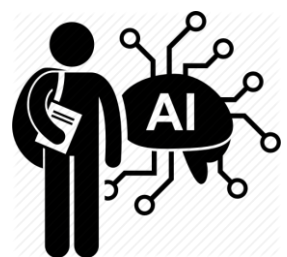
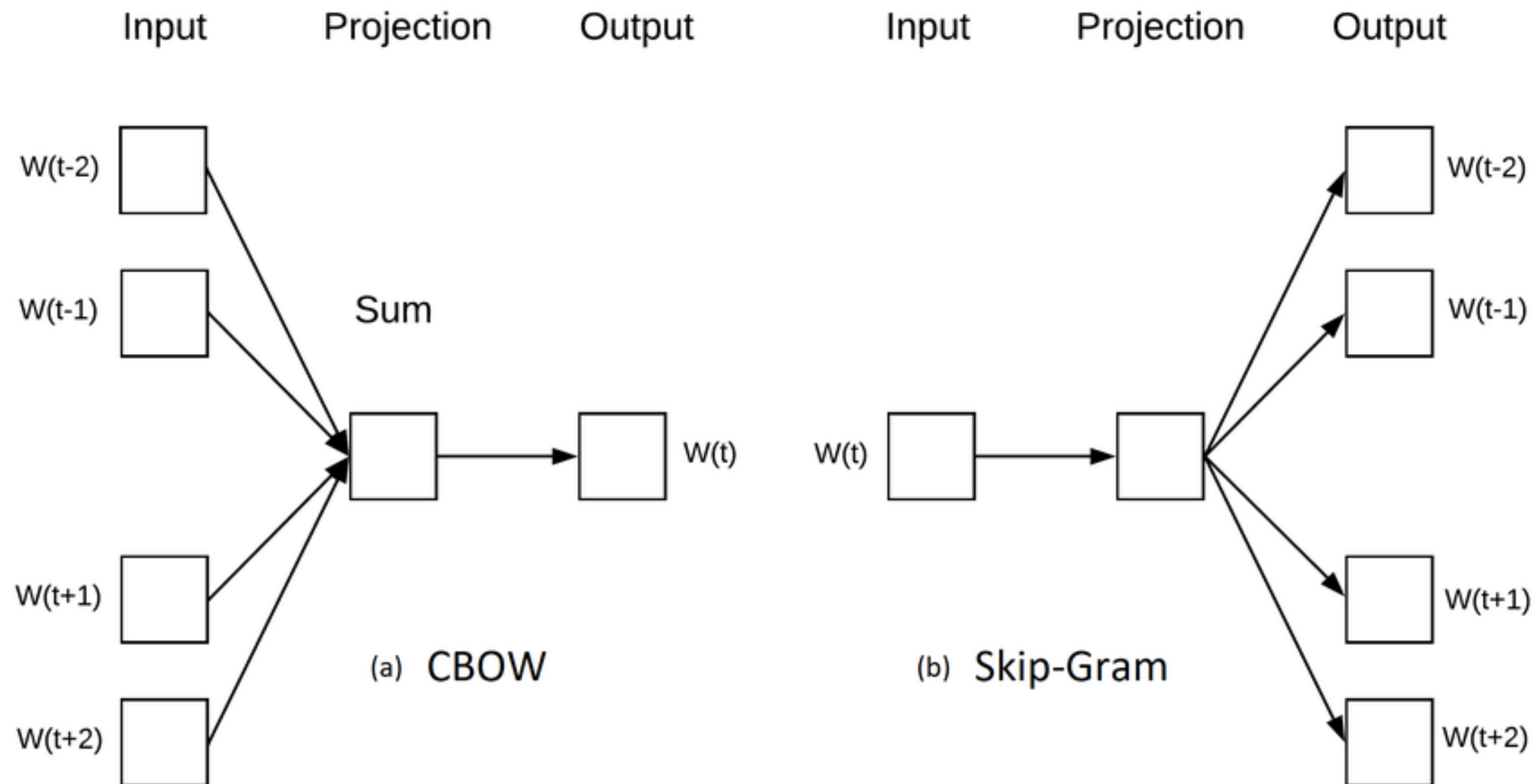


Source: <https://cambridge.org>

AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP

Word Embeddings

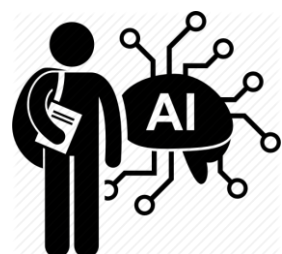
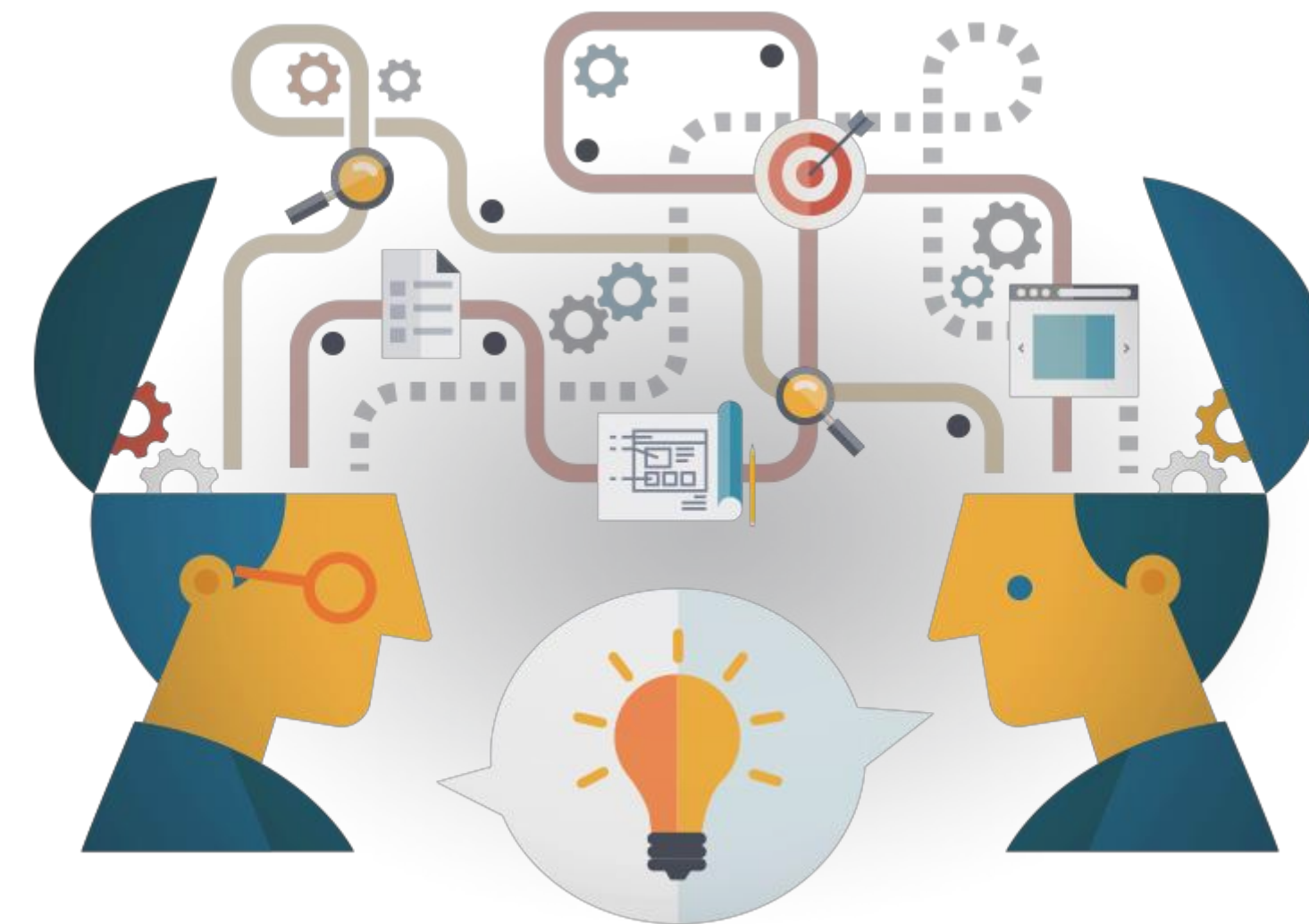
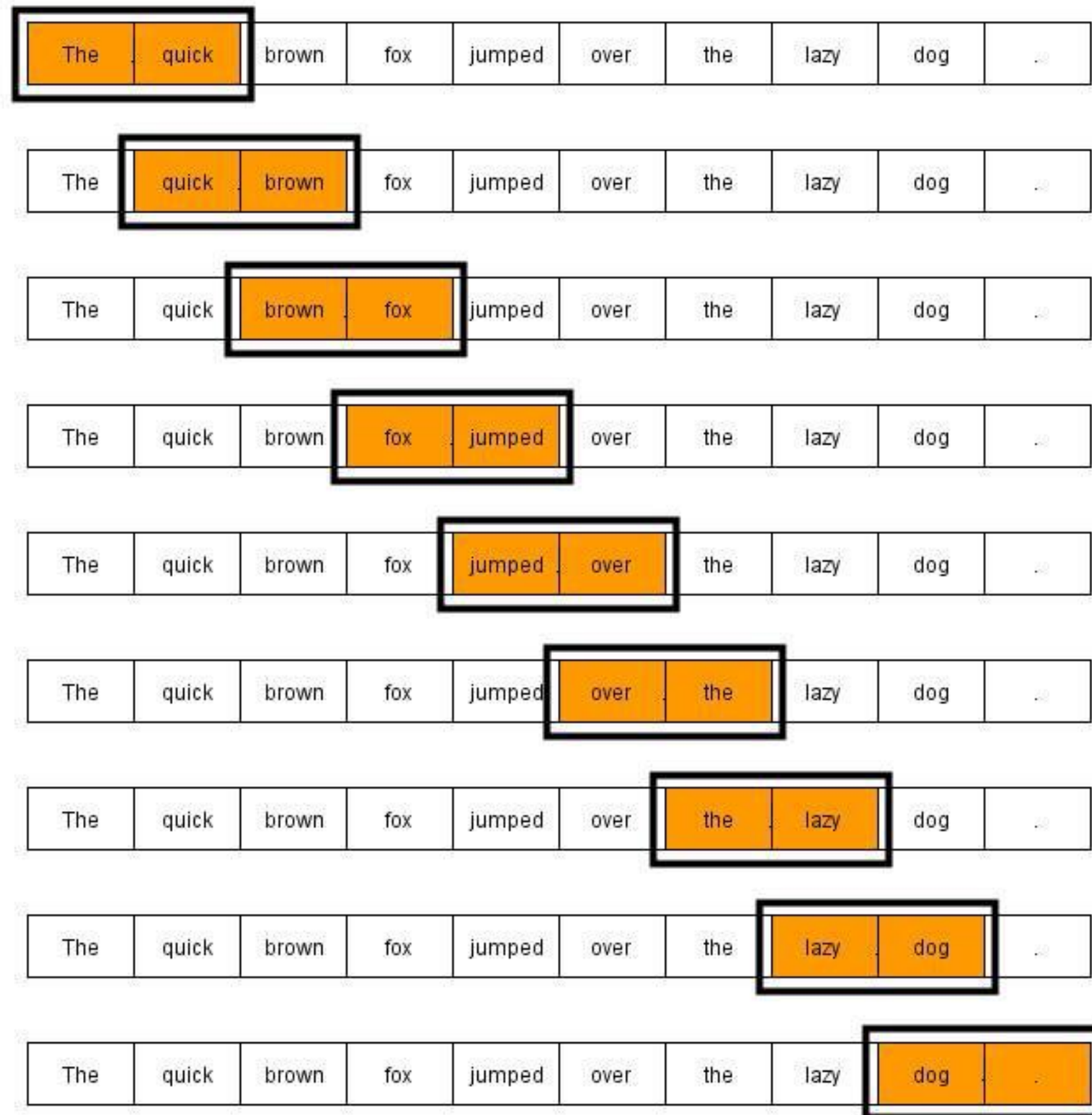


Source: <https://researchgate.org>

AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP

Transfer learning



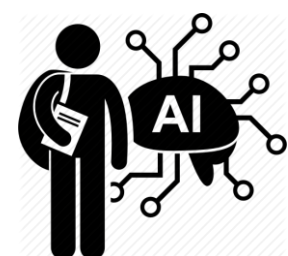
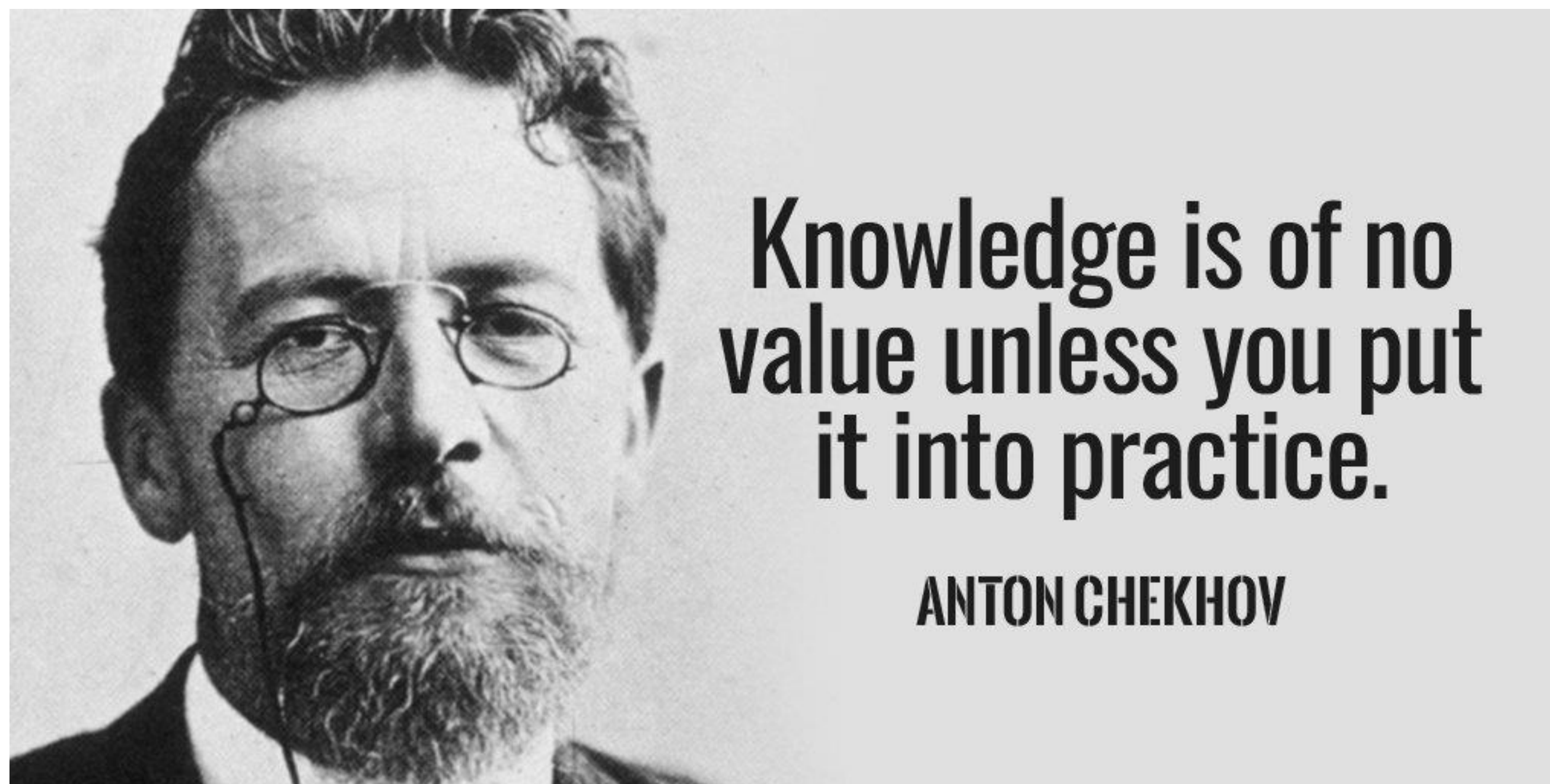
Source: <https://towardsdatascience.com> and <https://medium.com/paper-club>

AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP

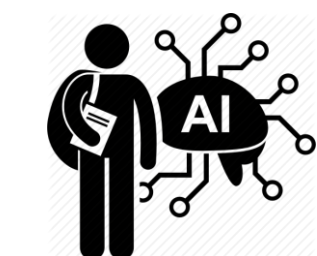


DEMO





Thank you



AI Day Wrocław
Volvo Group IT in Poland
29th November 2018

VOLVO
VOLVO GROUP